

Clustering

Cluster Analysis

- **Partitioning** a set of data objects into subsets or **clusters**
 - objects in a cluster are similar, yet dissimilar to objects in other clusters

انا معرفش ايه cluster اللي هاتطلع معايا فبحاول اقسام الداتا لمجموعه من cluster بحيث الداتا في cluster الواحد بتكون شبه بعض clusters ذات نفسها بتكون مختلفه. لو انا في شركه وعندي عملاء وبيشتروا منتجات فانا عايز اعرف مين اليوزر اللي شبه بعض عشان احصلهم اعلاناتاو منتجات معينه جنب بعض عشان اسهل عليهم.

- **Goal:** discovery of **previously unknown groups within the data**

Cluster بيستكشف حاجات نسبتها قليله في الداتا

(لو عندي مرضي بيعانوا من اعراض غير طبيعيه)

- Clusters are **implicit classes**
- **Applications** → business intelligence, image pattern recognition, web search, biology, security

ممكن استخدم cluster في السيكيورتي لو جالي اكثر من ريكوست علي السيرفر فهاقسم الريكوستات دي الي مين اللي عايز خدمات ومين اللي عايز يعمل هاك علي السيرفر

- Clustering can be used for *preprocessing* and *outlier detection*

ممكن استخدم cluster في preprocessing وتحديد outlier

Requirements for Cluster Analysis

- **Scalability** → currently handles small datasets, uses sampling

يعني complexity بيتزيد linear مع زياده الداتا مش موجود عن cluster algorithm كثير

- **Handling different attribute types** → mostly numerical

اتعامل مع انواع attribute المختلفه (اغلب الانواع numerical)

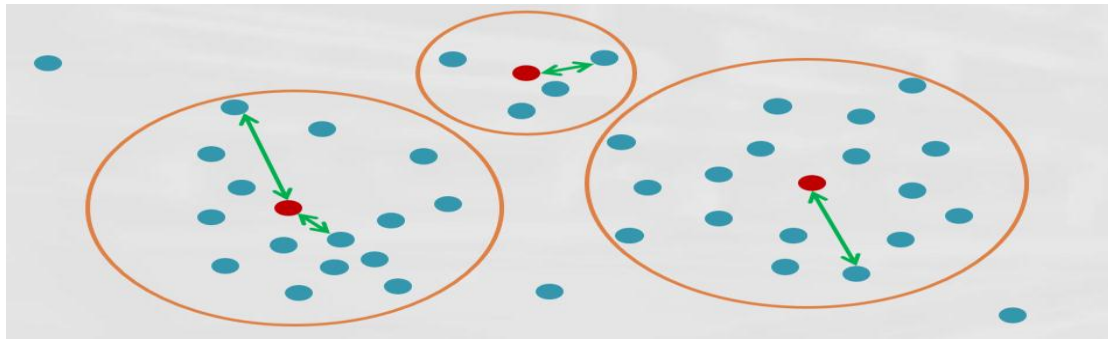
- **Discovering clusters with arbitrary shape** → currently mostly spherical

لازم اقدر اكتشف cluster باشكال مختلفه

- **Domain knowledge & input parameters** → # clusters & clustering results

لازم يبيقي عندي Domain knowledge مثلا الجوريزم k-means انا اللي بدخله k فاكيد هاتأثر علي الداتا. معظم الالجوريزمات بتعتمد علي parameter فلازم اعمل test عليهم

- **Handling noisy data** → currently sensitive to noise



اي noise ممكن تاثير علي cluster shape

لو بعمل cluster لليوزر في real time مش هاقدر اتعامل معاها كويس

- **Incremental clustering & insensitivity to input order** → new data requires re-computing clusters from scratch – sensitive to order

ترتيب الداتا لو اختلف ترتيب record هاختلف الناتج معايا

- **Handling high-dimensionality data** → mostly low Dimensionality

Cluster يقدر ي handle الداتا الي 3 dimension

- **Constraint-based clustering** → little support for domain constraints

معرفش احط exception لما بحط constrain لل cluster

- **Interpretability & usability** → are results comprehensible & usable?

ممكن يطلع cluster بس مش مفهومه ومقدرش استخدمها

Comparing Cluster Analysis Methods

- **The partitioning criteria** – *flat* or *hierarchical*?

هاقسم الداتا ل (flat) يعني one layer ولا الي (hierarchical) يعني اكثر من layer

- **Separation of clusters** – *mutually exclusive* or *overlapping*?

تقسيم ال cluster

mutually exclusive: يعني اقسام ال cluster بحيث ان لو طلعت cluster معين مايطلعش معايا ثاني في التقسيم

Overlapping: لو ممكن يكرر

- **Similarity measure** – *distance* or *connectivity/density*?

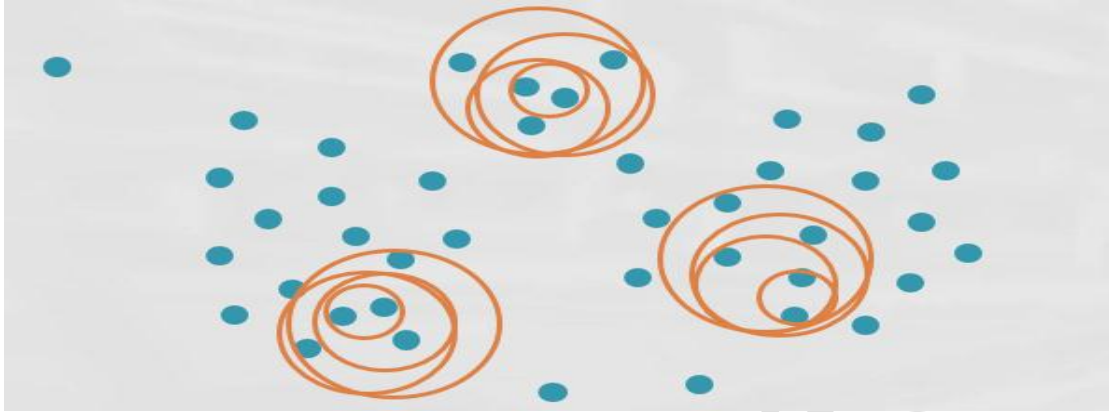
هايحسب similarity

Overview of Cluster Analysis Methods

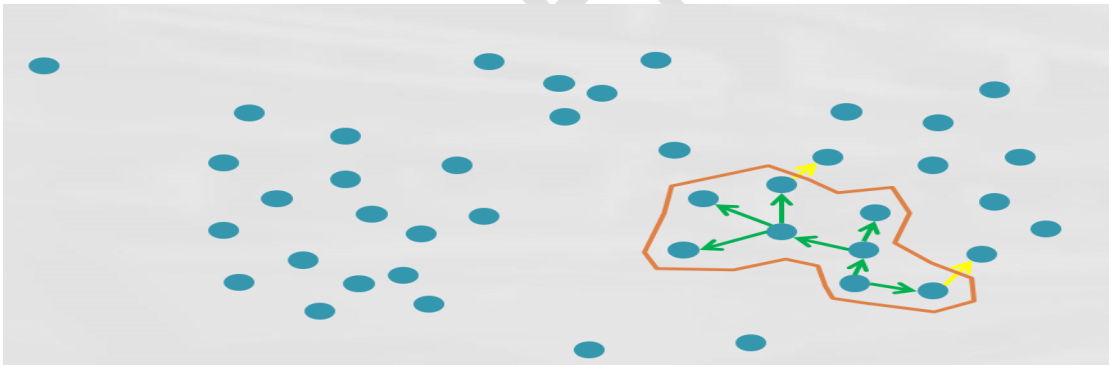
1. Partitioning

الاول هايغترض عدد cluster اللي عايزهم ويفرض مثلا ٣ نقط كل منهم في cluster
ويحسب بناء عليهم كل cluster والنقط اللي فيه زي k-means كده
ال similarity في cluster الواحد قريبين من بعض

2. Hierarchical



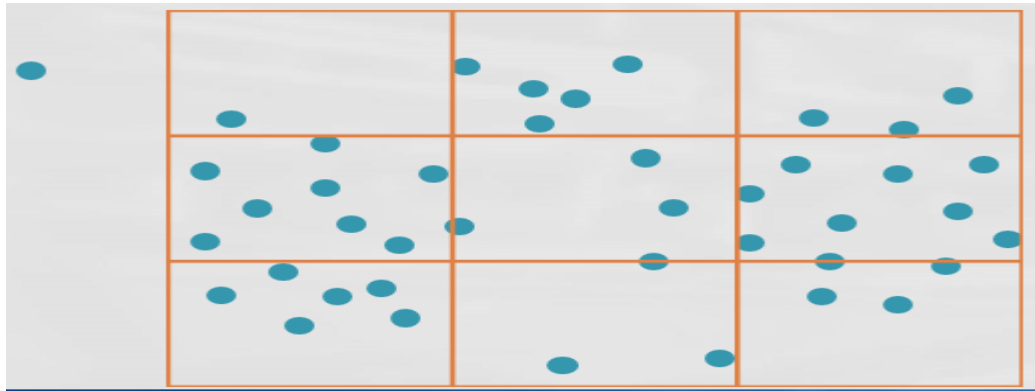
هافترض ان كل item هو cluster لوحده وهابدا ادمج كل اثنين مع بعض بحيث ان
ال اثنين يكونوا قريبين في similarity وهكذا لحد ماخلص كل النقط
مشكله Hierarchical لو طلعت مستوي معرفش انزل منه



3. Density-based

وده بيستخدم عشان يجيب الاشكال الغير منتظمه (non linear) ومن امثله DBSCAN
الاول بيبدأ باي نقطه وبغدين يبدأ احسب المسافه بينها وبين كل اللي حوالها علي حسب
المسافه اللي محددها واقل عدد نقط موجود في الكلستر الواحد بتبقى متحده قبل مااشتغل
ويحسب المسافه بين كل نقطه واللي حوالها لو اللي اشوف هل كل النقط اللي حوالها واقل
من المسافه اللي متحده هل هما نفس العدد يعني لو انا محدد المسافه 2 واقل عدد في
الكلستر الواحد 4 يبقى اشوف المسافه بين النقطه واللي حوالها لو لقيتهم اقل من المسافه
اللي محددها وكمان يكون عددهم اكبر من او يساوي اقل عدد انا محدده

Grid-based



ب map ال grid علي الداتا اللي عندي ويعمل fixed cluster وبخلي كل cell لوحدها

Overview of Cluster Analysis Methods

Method	Characteristics
Partitioning methods	<ul style="list-style-type: none"> — Find <u>mutually exclusive</u> clusters of <u>spherical shape</u> — <u>Distance-based</u> — May <u>use mean or medoid</u> to represent cluster center — Effective for <u>small- to medium-size data sets</u>
Hierarchical methods	<ul style="list-style-type: none"> — Clustering is <u>hierarchy</u> involving multiple levels — Cannot correct <u>erroneous merges/splits</u> — May consider object "<u>linkages</u>"
Density-based methods	<ul style="list-style-type: none"> — Can find <u>arbitrarily shaped clusters</u> — Clusters are <u>dense regions</u> separated by <u>low-density regions</u> — Each point must have a <u>minimum number of points within its "neighborhood"</u> — May <u>filter out outliers</u>
Grid-based methods	<ul style="list-style-type: none"> — Use a multi-resolution <u>grid data structure</u> — <u>Fast processing time</u>

Partitioning Methods

1. K-Means – A Centroid-Based Technique

- Divide dataset into k **mutually exclusive** clusters

بيقسم الداتا لمجموعه من الكلستر ويكونوا مختلفين عن بعض وانا اللي بحدد عدد cluster

- Clusters are represented by their **centroids**

كل cluster بيتمثل بالسنتر بتاعه

- A centroid is a **cluster's center point**

- In k -means \rightarrow centroid is **mean** of points within cluster

في k -mean السنتر بيتمثل بالنقط اللي داخل cluster

- Each object x in cluster has a distance from centroid $c_i \rightarrow dist(x, c_i)$
- x is assigned to most similar cluster $\rightarrow C_i$ with **min** $dist(x, c_i)$
- Cluster means are updated, then assignment is repeated
- To measure cluster quality \rightarrow minimize sum of squared errors

$$E = \sum_{i=1}^k \sum_{x \in C_i} dist(x, c_i)^2$$

Factors to consider:

- Selection of k
- Selection of **initial centroids**
- Calculation of **dissimilarity**
- Calculation of **cluster means**

When it fails!

العوامل المؤثرة في الالجرويزم ده

اختياري K و انا اللي بافترض السينتر بتاع كل cluster وحساب **dissimilarity** و **cluster means**

- Clusters with **very** different sizes & with **concave** shapes



مايطلعش cluster مختلفين في الحجم او شكلهم مقعر او مش مستوي

Hierarchical Methods

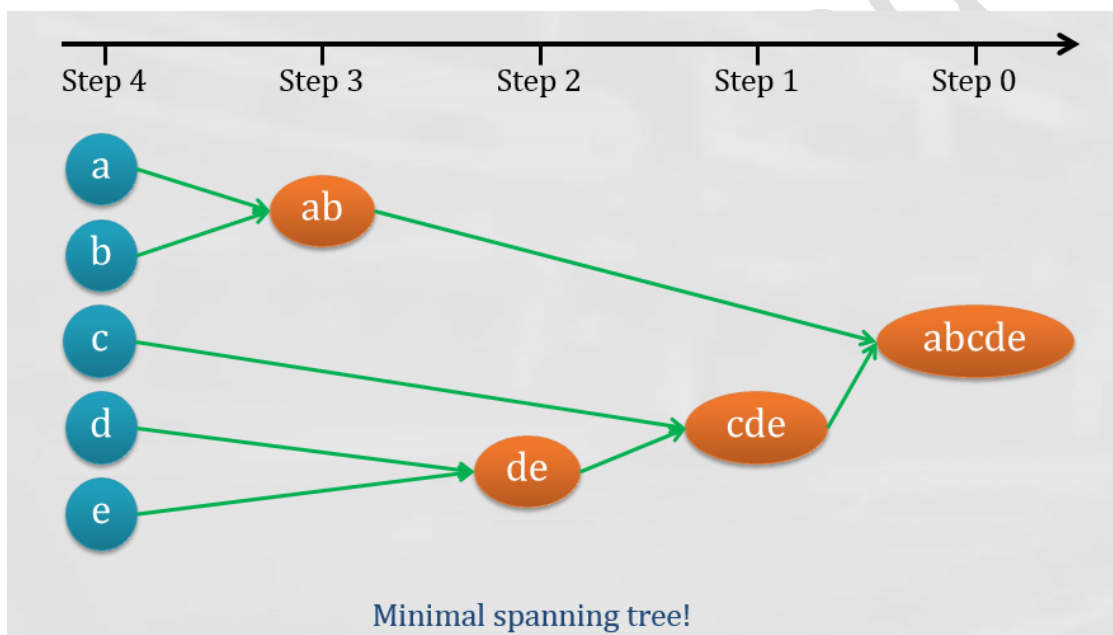
1. Agglomerative versus Divisive Clustering

• **Hierarchical clustering** → group data objects into a *hierarchy* or “tree” of clusters

• **Agglomerative** → bottom-up (merge) composition

- Each object has its own cluster
- Two clusters that are closest **merged** into a bigger cluster
- Iteratively merge till *termination condition* or *single cluster* is formed

بافتراض ان كل item او object وكل اثنين cluster قريبين من بعض يبقوا cluster واحد وافضل اجمع كل اثنين cluster لحد شرط معين اكون فرضته في الاول او عدد cluster معين

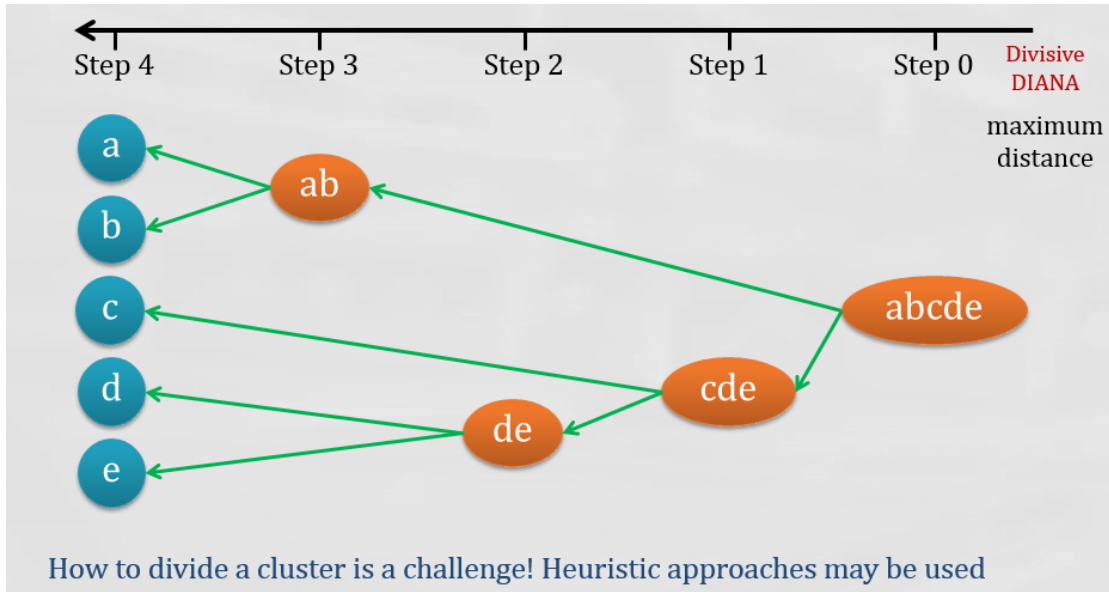


مثلا لو قيست المسافه بين كل نقطه والثانيه ولقيت ان a , b الاقرب لبعض وكذلك c و d, e دلوقتى انا عندي ٣ cluster هاعمل نفس حكاية بينهم هاشوف المسافه بين c و ab والمسافه بين c و de لما اقيس المسافه هاخذ النقطتين بتوع ab واقيس كل واحده مع c وكذلك de وهاشوف في الاخر c اقرب لاي نقطه يبقى تبع cluster اللي فيه النقطه دي وهكذا

• **Divisive** → top-down (split) composition

- All objects in one big cluster
- **Divide** into subclusters
- Recursively divide subclusters into even smaller subclusters
- Terminate when *each object has his own cluster* or *objects in clusters are similar “enough”*

هنا العكس هافترض ان كل الدانا اللي عندي عبارة عن cluster واحد وابدا اقسام فيها لحد



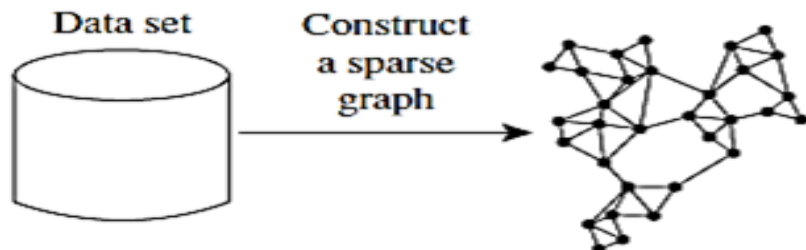
شرط معين او لما الاقي كل object في cluster متشابهين
المشكلة هنا هاقسم بناء علي ايه؟؟؟

CHAMELEON: Multiphase Hierarchical Clustering Using Dynamic Modeling

- Cluster similarity based on: مبني علي
 - Interconnectivity** → how well connected objects are within a cluster
مدي ان object متوصله ببعض داخل cluster الواحد
 - Closeness** → the proximity of clusters
مقدار التقارب بين cluster

1. Construct a **sparse graph**:

- Vertices are data objects
- there exists an edge between two vertices $\{x, y\}$ if x is among the k -most similar objects to y → k -nearest neighbor graph



- Edges are *weighted* to reflect similarity

كل edge عليها weight بتاعها اللي بيمثل المساحه بين 2 node وارسم graph بناء علي قواعد معينه بيستخدم k -nearest neighbor graph

2. Graph partitioning:

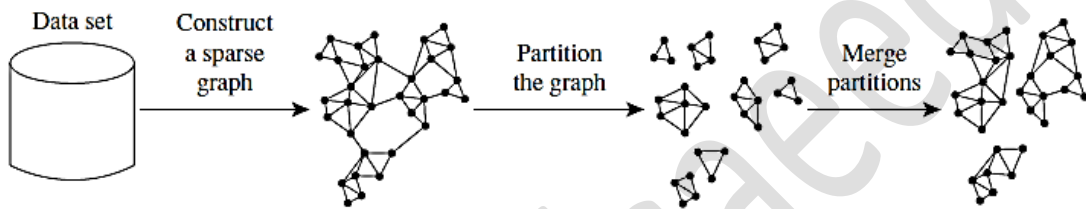
- Minimize **edge cut** $EC(C_i, C_j) \rightarrow$ minimize weight of edges that would be cut to split C into C_i & $C_j \rightarrow$ measures *absolute interconnectivity*

عايز اقسام graph عشان اقلل نسبه الخطا عندي

Width عايزه اقل مايمكن فيه الجوريزمات بتعمل الكلام ده

3. Agglomerative clustering:

- Measure *relative interconnectivity* $RI(C_i, C_j)$
- Measure *relative closeness* $RC(C_i, C_j)$



بعمل graph للداتا اللي عندي وبعد كده بقسم graph ده لاشكال اصغر منه عشان اقلل نسبه الخطا وبعد كده ادمجهم ببعض بحيث اطلع cluster من كل كام graph صغير Graph الاول ده مش انا اللي بعمله في applications جاهزه بتعمل الكلام ده.

- Relative Interconnectivity** $RI(C_i, C_j) \rightarrow$ absolute connectivity between C_i, C_j normalized by *internal interconnectivity* of C_i, C_j

$$RI(C_i, C_j) = \frac{|EC_{\{C_i, C_j\}}|}{\frac{1}{2}(|EC_{C_i}| + |EC_{C_j}|)}$$

where $EC_{\{C_i, C_j\}}$ = sum of weights of edges that connect C_i with C_j .

- $EC_{C_i} \rightarrow$ min sum of cut edges that partition C_i into two roughly equal parts

Relative Interconnectivity->Sum

- Relative Closeness** $RC(C_i, C_j) \rightarrow$ absolute closeness between C_i, C_j normalized by the internal closeness of C_i, C_j

$$RC(C_i, C_j) = \frac{\overline{S_{EC_{\{C_i, C_j\}}}}}{\frac{|C_i|}{|C_i| + |C_j|} \overline{S_{EC_{C_i}}} + \frac{|C_j|}{|C_i| + |C_j|} \overline{S_{EC_{C_j}}}}$$

- $\overline{S_{EC_{\{C_i, C_j\}}}}$ → average weight of edges connecting vertices in C_i to vertices in C_j

هو $\overline{S_{EC_{\{C_i, C_j\}}}}$ average weight لل edge اللي بيربط مابين

Relative Closeness->average

- $\overline{S_{EC_{C_i}}}$ → avg weight of edges belonging to **min-cut bisector** of C_i

شرح المعادلات

$\overline{S_{EC_{\{C_i, C_j\}}}}$ and $EC_{\{C_i, C_j\}}$ يبقوا edges اللي بتربط بين 2 كلستر

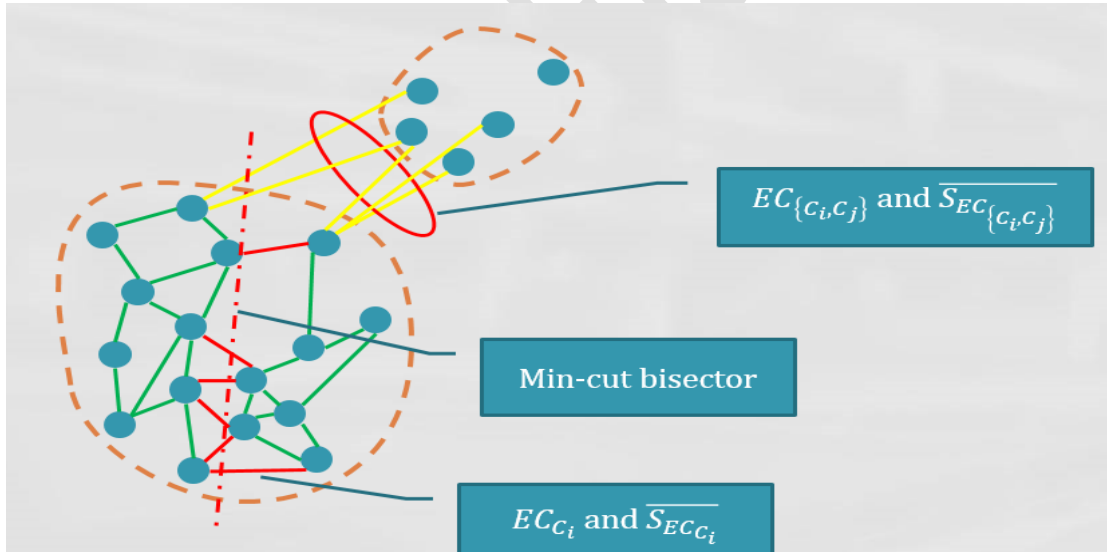
$|EC_{C_i}| + |EC_{C_j}|$ or $\overline{S_{EC_{C_i}}} + \overline{S_{EC_{C_j}}}$ = weighted sum of edges that partition the cluster into roughly equal parts

ودول ببقوا edge اللي جواه cluster اللي بيقسموه الي جزئين متساوين

من الرسمه **min-cut bisector** بيقسم cluster اللي جزئين ال edge اللي بيمر بيها

الخط ده هما دول $|EC_{C_i}| + |EC_{C_j}|$ or $\overline{S_{EC_{C_i}}} + \overline{S_{EC_{C_j}}}$

خلاصه المعادلات ان اللي في البسط هو مابين cluster والمقام كل cluster لوحدها



Density-based Methods

DBSCAN: Density-Based Clustering Based on Connected Regions with High Density

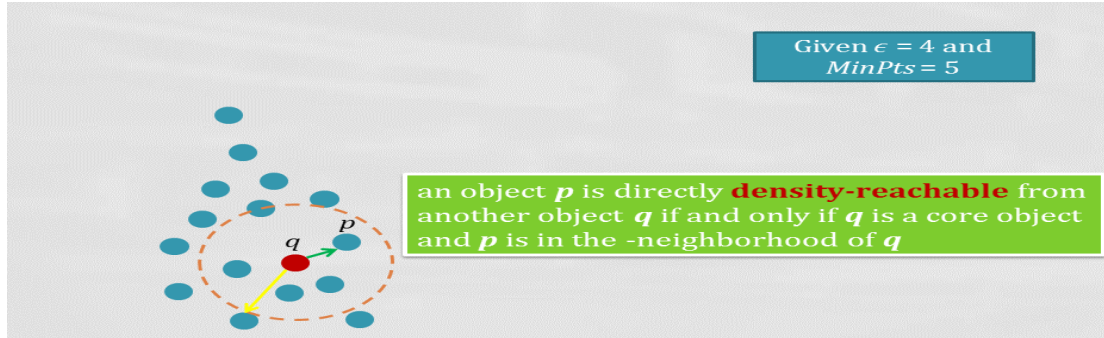
- Find core objects (with dense neighbourhoods)
- Connect core objects to form dense clusters
- User provides:
 - ϵ -neighborhood of object o → space within a radius ϵ centered at o المسافه (نصف القطر) من النقطه لاي نقطه حوالها

● **Neighborhood density** → # objects in that neighborhood

● **MinPts** → density threshold for a neighborhood في عدد من النقط في الكلستر الواحد

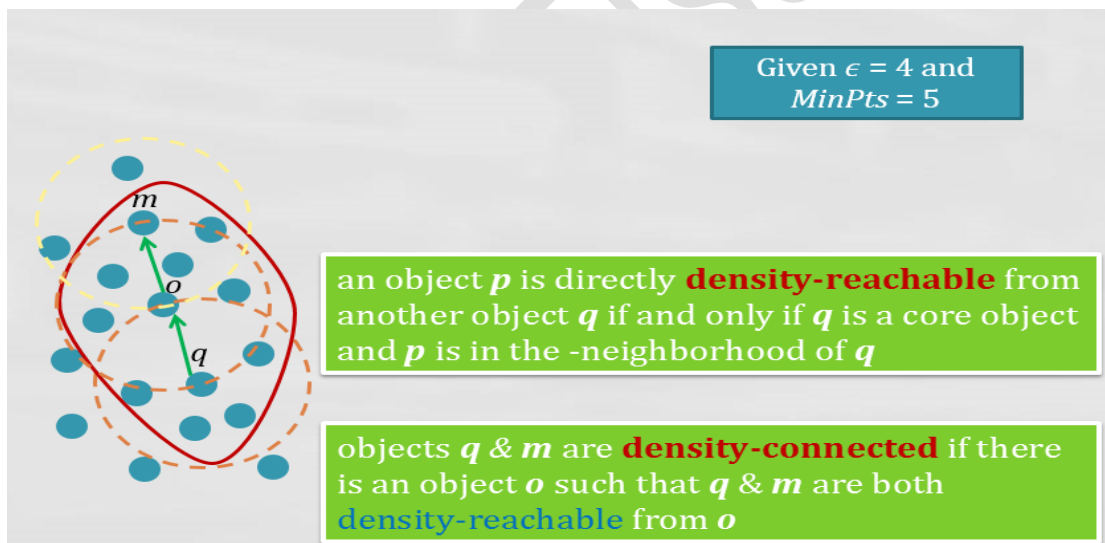
● **Core object** → object whose ϵ -neighborhood contains at least *MinPts* object

النقطة التي اقل عدد من النقط حوالها في دائره نصف قطرها المسافه التي محددها هو اكبر من او يساوي *MinPts*



p is **density-reachable** from *q* (يعني *p* اقدر اوصلها من *q*)

يعني لو انا محدد المسافه 4 (نص القطر) واقل عدد في الكلستر الواحد 5 يبقي اشوف المسافه بين النقطه واللي حوالها (اللي حوالها في دائره نصف قطرها المسافه التي



محددها وهي 4) لو في الدائره دي عدد النقط اقل من 5 يبقي اشوف نقطه ثانيه وتبقي دي السينتر بتاع الكلستر واحد من خلالها الكلستر لو فيه نقطه علي خط الدائره تتحسب انها تبعها

هاشوف كل النقط واطلع clusters اللي موجوده

q & *m* **density-connected** *o* لو المسافه مابين *q* و *o* زي المسافه بين *m* و *o* يعني اقدر اوصل لل *q* من *o* او اقدر اوصل لل *m* من *o*

مميزاته بيقدّر يطلع cluster باشكال مختلفه

Evaluation of Clustering

Assessing Clustering Tendency

- Determines whether a given data set has a non-random structure
- **Hopkins Statistic** → Statistical tests for spatial randomness

- Sample n points, p_1, \dots, p_n uniformly from D
- For each point, p_i find its nearest neighbor in D

بيقارن كل dataset باقرب النقط ليها (neighbors) ويحسب المسافه بينهم ويجيب اقربهم. لو عندي n فيها شويه نقط p_i بيشوف لكل نقطه فيهم اقرب النقط في dataset

$$distance\ x_i = \min_{v \in D} \{dist(p_i, v)\}$$

- Sample n points, q_1, \dots, q_n , uniformly from D
- For each q_i find its nearest neighbor of q_i in $D - \{q_i\}$

بيقارن كل dataset باقرب النقط ليها (neighbors) ويحسب المسافه بينهم ويجيب اقربهم. لو عندي n فيها شويه نقط p_i بيشوف لكل نقطه فيهم اقرب النقط في $D - \{q_i\}$

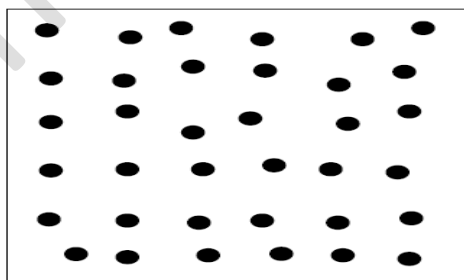
$$distance\ y_i = \min_{v \in D, v \neq q_i} \{dist(q_i, v)\}$$

- Calculate the Hopkins Statistic H :

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

- If D is uniformly distributed, $\sum_{i=1}^n x_i$ and $\sum_{i=1}^n y_i$ are roughly equal, and $H \approx 0.5$

وبيكرر العمليه دي اكر من مرة عشان يقدر يحسب اقرب النقط ليها



Measuring Clustering Quality – Extrinsic Methods

Extrinsic methods → compare clustering against ground truth (supervision)

النوع ده ببيقي supervised ويتم عن طريق مقارنه الكلستر ب ground truth

- Assign a score $Q(C, C_g)$ to capture:
 - **Cluster homogeneity** → the purer the better – clusters represent separate class labels

تجانس الكليستر كل لما يكون انقي ومتجانس مع بعضه كل لما يكون افضل

- **Cluster completeness** → an object with a class label belongs to the cluster representing that class label

يكون object في class يعبر او يمثل ال class كله

يعني كل class يتمثل في مجموعه قليله من clusters

- **Rag bag** → objects that can't be merged into clusters belong to a *rag bag* – penalize a *misc. object* when put in a *pure cluster* more than in a *rag bag*

Object اللي ماينفعش احطها في cluster بتنتمي ل Rag bag

- **Small cluster preservation** → splitting a small category is *more harmful* than splitting a large category

اني اقسام category صغير الي اجزاء اصغر مضر اكثر من اني اقسام category كبيره

- Ex. **Bcubed precision and recall** of every object in dataset:

- Precision → how many objects in the same cluster ∈ the same category as the object

كام عدد object في cluster الواحد ينتموا لنفس category

- Recall → how many objects of the same category are assigned to the same cluster

كام عدد object في category الواحد ينتموا لنفس cluster

Intrinsic methods → measure how well the clusters are separated

النوع ده بيبقي unsupervised و ground truth غير متاح هنا بيقوم جوده الكليستر بانه بيشوف ازاي clusters منفصلين وازاي بيتم الدمج بينهم

- Ex. **The silhouette coefficient** → difference between:

- *average distance* between object **o** and all other objects in the cluster to which **o** belongs (captures cluster correctness) – smaller is better (more compact)

متوسط المسافه بين object O وبين كل object اللي في نفس cluster اللي بيتنتمي ليها O كل لما كانت اصغر كل لما كانت افضل (العلاقه بين object في نفس الكلاس)

- *minimum average distance* from **o** to all clusters to which **o** does not belong (captures degree of separation from other clusters) – larger is better

اقل متوسط مسافه بين object O وبين كل cluster التانيه اللي مش فيها O كل لما كانت اكبر كل لما اكنت احسن (العلاقه بين object في كلاس وبقية الكلاسات)

- Compute average silhouette coefficient for all objects in a cluster **or** over all of the dataset

- **+ve** → clustering is good

- **-ve** → clustering is bad

بحسب المتوسط لل silhouette coefficient لكل dataset او لكل عنصر في cluster
واشوفه لو موجب يبقى عمليه cluster كويسه لو سالب تبقي فيه مشكله في عمليه
cluster

Ahmad ElSaeed